

Problem 1.

Consider a system with $\gamma=0.9$ the following state and action spaces:
 $S = \{-1, 1, 2\}$, $A = \{-1, 0, 1\}$. The available batch data are as follows:

$$D = \{ (s_0 = 1, a_0 = 1, r_1 = 1, s_1 = 2), (s_1 = 2, a_1 = 0, r_2 = -1, s_2 = 1), \\ (s_2 = 1, a_2 = -1, r_3 = 0, s_3 = -1) \}$$

consider the basis function $\Phi(s, a) = a^2 s + a s \rightarrow a$ with initial weights $w^0 = 1$.
 Perform LSPI to compute w' and w'' , and policy associated to w'' .

* In case of tie for action selection, give the preference to -1, then 0 and finally 1.

Example $\rightarrow \arg \max_{a \in \{-1, 0, 1\}} \left\{ \frac{-1}{2}, \frac{0}{2}, \frac{+1}{2} \right\} = -1$

Problem 2.

Repeat Problem 1 using the basis function $\Phi(s, a) = \begin{bmatrix} a s \rightarrow a \\ a^2 s \end{bmatrix}$ with initial weights $w^0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$. Perform LSPI to compute w' and w'' , and policy associated

to w'' . Is the final policy (i.e., π^2) different from problem 1?

Can two basis functions, in general, lead to different policies?

Appendix: LSPI.

Policy Evaluation

$$\omega_{k+1}^- \rightarrow Q_{(s,a)}^+ = \Phi_{(s,a)}^T \omega^-$$

$$\pi_{(s)}^- = \operatorname{argmax}_{a \in A} Q_{(s,a)}^+ = \operatorname{argmax}_{a \in A} \Phi_{(s,a)}^T \omega^-$$

$$\omega^- = \omega^+$$

Policy Improvement

$$A = \frac{1}{L} \sum_{i=1}^L \Phi_{(s_i, a_i)} [\Phi_{(s_i, a_i)} - \gamma \Phi_{(s_{i+1}, \pi_{(s_{i+1})}^-)}]^T$$

$$b = \frac{1}{L} \sum_{i=1}^L \Phi_{(s_i, a_i)} r_{i+1}$$

$$\omega^+ = A^{-1} b$$